

NON-UNIQUENESS OF ATMOSPHERIC MODELING

PHILIP G. JUDGE and SCOTT W. MCINTOSH*

*High Altitude Observatory, National Center for Atmospheric Research[†], P.O. Box 3000,
Boulder CO 80307-3000, USA*

(Accepted 18 October 1999)

Abstract. We focus on the deceptively simple question: how can we use the emitted photons to extract meaningful information on the transition region and corona? Using examples, we conclude that the only safe way to proceed is through forward models. In this way, inherent non-uniqueness is handled by adding information through explicit physical assumptions and restrictions made in the modeling procedure. The alternative, ‘inverse’ approaches, including (as a restricted subset) many standard ‘spectral diagnostic techniques’, rely on more subjective choices that have, as yet, no clear theoretical support. Emphasis is on the solar transition region, but necessarily discussing the corona, and with implications for more general problems concerning the use of photons to diagnose plasma conditions.

As in all other astronomical endeavors where the subject is not directly retrievable, the study of the transition region ranges from purely *ab-initio* physical models to semi-empirical conclusions drawn more directly from observational data. I have used the expression ‘semi-empirical’ because *any* conclusions concerning the state of the transition region material must rely on physical models and implicit assumptions. Only the photons are directly observable.

L. S. Anderson-Huang (1998)

1. Introduction

In this exciting era of multiple space experiments (TRACE, SOHO) devoted to studying photon spectra emitted by the Sun, it is appropriate to revisit the important question: *How can we use the emitted photons, with simplified models, to extract objectively meaningful information about the emitting plasmas?* The purpose of this review is to look into this question, with particular emphasis on one of the observationally best-studied, but least understood, regions of the Sun’s atmosphere: the transition region (henceforth ‘TR’). By using the term ‘the solar transition region’, we mean all plasma that contributes significantly to radiation from ions whose emission, under coronal ionization equilibrium conditions, peaks between say 2×10^4 and 5×10^5 K. We will not deal explicitly with flares or ‘explosive events’. A review of the TR, containing some aspects of the issues discussed here, is given by Anderson-Huang (1998).

[†]The National Center for Atmospheric Research is sponsored by the National Science Foundation.

*Also with the Advanced Study Program of NCAR.



The article is organized as follows. In Section 2 we briefly review the evolution of ideas on the structure of the solar TR. Section 3 discusses approaches one can take towards inferring properties of plasmas from their emitted radiation. Section 4 discusses specific examples of applications of the various approaches, by using the ‘forward-inverse’ approach in which simulated data are subjected to some traditional ‘spectral diagnostic techniques’. The paper closes with a review of lessons learned from these examples, and discusses the merits of the various approaches. It is argued that the only safe way to progress is through forward models.

2. Current Pictures of the Solar Transition Region

From a physical point of view, one would like to be able to answer the following questions. What is the nature of the sources of mass, momentum and energy for the chromosphere/TR/corona? How does the chromosphere/TR/corona respond? Answers are not yet forthcoming because of limitations in both observations and theory (e.g., the physics of reconnection is an active research area). We must therefore seek answers to more restricted questions, using both observations and simple physical models to make progress. One such question is, quite simply, what is the basic structure of the TR?

Two classes of models are currently considered to be important for describing the solar TR: classical TR models (‘CTR’), in which the emission is formed at the thermal interface between the chromosphere and corona, and other non-CTR (‘nCTR’) models, in which the emission from TR ions comes from an entirely different structure. There is (perhaps surprisingly) active debate concerning the relative contributions of these different pictures to the observed TR emission. One aim of this paper is to examine critically the arguments for and against each class of model. We will try to show that both models have problems and merits. We will also argue that one should avoid prejudices when analyzing solar data, since the information in the observations alone is not enough to discriminate between the models, and current theory is not yet able to provide answers.

CTR models are represented in the classic work of Gabriel (1976). Variations on simple CTR models prompted by their (well-known) failure to produce enough radiation emitted below 10^5 K have been presented by Athay (1990), Cally (1990) and Ji, Song, and Hu (1996). Evidence for the failure of CTR models, and for support of non-CTR models is given by Feldman (1983, 1987), Dere *et al.* (1987), Feldman and Laming (1994), and Feldman (1998), among others. Physical non-CTR models, at least partially inspired by the above work, have been considered by Antiochos (1984), Rabin and Moore (1984), Antiochos and Noci (1986), Sturrock *et al.* (1990), Cally and Robb (1991), Roumeliotis (1991), and Spadaro, Lanza, and Antiochos (1996). These models are strikingly different. Consider, for example, the ‘thread-like’ structure envisaged by Dere *et al.* (1987), which has extremely small area filling factors, severely limiting the supply of mass into the corona through the

observed structures, whereas the model of Gabriel (1976) has an area filling factor closer to unity. It is thus important to review the evidence and arguments in support of these different physical pictures.

3. Approaches

Remotely sensed data are generally interpreted using variants of two different approaches: *forward* and *inverse* approaches (e.g., Craig and Brown, 1986). There is another popular approach which, as we will see, is a restricted form of the inverse method. This approach amounts almost to ‘common sense’ or ‘intuition’, but has (to our knowledge) never been properly defined. For want of a better name we will call these ‘empirical’ approaches (Anderson-Huang, 1998, prefers ‘semi-empirical’), and will try to define more exactly what these are, in terms of inverse methods. Figure 1 summarizes schematically the ideas behind the various approaches. The examples in Section 4 should serve to illustrate these concepts more clearly.

Forward methods are conceptually the simplest: one develops a physical model for the gas/plasma that emits the photons, compares with observations, and proceeds with modifications, or stops, based upon similarities or differences.

Inverse methods aim to determine solutions, or ranges of solutions, together with uncertainties, by applying a formal ‘inversion’ of the forward problem (Figure 1). Thus, starting with the observed quantities, the model (assumed known) is used to infer the ‘source function’.

Empirical methods are the most commonly used. Often they are applied without due awareness of the true underlying assumptions (present authors included, see e.g., Brage, Judge, and Brekke, 1996), and we can do no better than give a simple example. Consider a volume of plasma that emits photons in lines whose emission coefficients depend only on the plasma electron temperature T and density n , as $G(T) \times n^2$. It is well known that the inverse methods should be used to determine limits on the form of the function $f(T) = \xi(T)$, the differential emission measure, from a set of such lines, with kernels $K(T) = G(T)$. Yet it is also possible to determine a ‘mean temperature’ from just a pair of suitable lines, by simply asking the question: ‘what is the single temperature that is compatible with the data’? As can be seen from Figure 1, this amounts to making the *assumption* that the source term $\xi(T)$ can be approximated by a Dirac- δ function.

While this example seems like a limited case, many of the empirical approaches fall into this class. Thus it is common to assign directly one number for a parameter f from one measurement (or combination of measurements) g . For example one often reads ‘the velocity [density, temperature, abundance of element X . . .] of the plasma is . . . because the measurements are . . .’.

These considerations are important, and not simply of academic interest for several reasons:

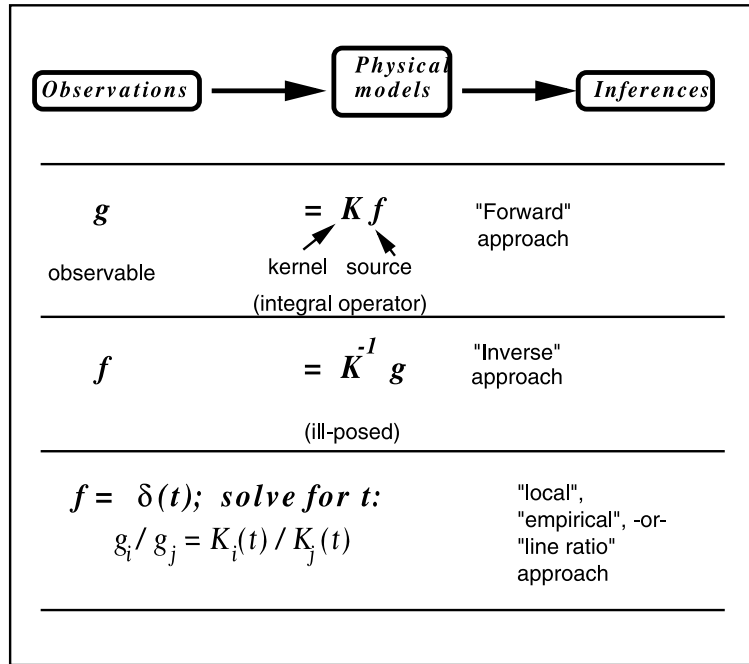


Figure 1. Approaches to the interpretation of remotely sensed data. g represents the observations (usually a data vector), K is an integral operator containing the physics of the emitting plasma (part of 'the model'), and f is the 'source term', some property of the emitting plasma that is desired (the other part of 'the model'). The model is here assumed to depend on independent variable t .

1. Photon spectra emitted from observable plasmas are integrals over volume because of radiative transfer. This applies even to astrophysical plasmas believed to be in some sense 'resolved'¹.
2. The integral operators (represented schematically by K in Figure 1) are such that photon spectra are (for related reasons) in general compatible with many different source terms. The inverse problem is often severely ill-posed.
3. All methods explicitly assume a model, they are all subject to non-uniqueness and ill-posedness, and all require additional information to constrain the solutions.

Points 1 and 2 show that *information must be added to the data to derive reliable information on the emitting plasma*. In particular, the intrinsic ill-posedness of the problem implies that one cannot assume that a one-to-one relationship exists between observations and properties of the emitting plasma, appealing though such a method is. We review the 'pros and cons' of these approaches.

The forward approach adds information through the physics put into the models, and the boundary conditions. An example might be a simulation of the gas

¹Note that even in the Sun, it is rare to find thermal linewidths for spectral lines, implying the presence of unresolved structure.

dynamics and radiation in magnetic flux tubes (e.g., Steiner *et al.*, 1998). Advantages are that this is physically based, and is generally applicable to within the known physics of the system. Disadvantages are that sometimes important physical processes are not understood (e.g., reconnection), boundary conditions may not be known (e.g., the nature of the photospheric ‘driver’, Parker, 1988), and understanding such processes might be the goal of the study! Furthermore forward calculations must often be over-simplified to the extent that they may be far from reality, for example by reducing an intrinsically 3D problem to 1 dimension.

The inverse approach assumes that the spectrum formation really ‘fits’ into the inverse formalism, and therefore that the inversion makes physical sense. The inverse solution (e.g., the differential emission measure $\xi(T)$ as determined from a set of emission line intensities) is often represented as the ‘Holy Grail’ of remotely sensed data (Craig and Brown, 1986), as indeed it is, provided the forward problem really is of the correct form². An advantage of inverse theory is that it yields constraints on the *range* of solutions compatible with the data. Disadvantages are that the emitting source may not comply with constraints needed for the spectrum formation to be written in suitable form (i.e., a Fredholm equation. The $\xi(T)$ inverse problem for the TR may not in fact be written in this form, Judge *et al.*, 1995). Further, one must add sometimes unphysical constraints to ‘regularize’ the solution (e.g., Craig and Brown, 1986), i.e. deal with the ill-posedness. One example is the use of restricted splines or low order derivatives to determine $\xi(T)$ (e.g., see the articles in Harrison and Thompson, 1991).

The empirical approach adds information through a radical assumption of the form of the solution – as we have seen the idea that ‘because of observable g then physical parameter X is. . .’ amounts to assuming that the solution is in fact a δ -function. Often one hears support for this method based upon Occam’s Razor. Advantages are that this method is easy, and popular, but there are severe disadvantages. Simply put, the very basis of the method involves a drastic simplification of a difficult and ill-posed inverse problem, this amounts (in essence) to subjective choice, and there is absolutely no measure of uniqueness in the interpretation. Another drawback of the approach is its ease of application, which naturally leads to popularity and (in our opinion) unfortunately some acceptance. As a community we should be aware that other interpretations are possible and that we might be guilty of a ‘collective mis-interpretation’ of data. Other plasma conditions might be equally compatible with the data.

We turn to some examples of interest to the Sun to try to solidify these ideas.

²Furthermore the ill-posedness means that the issue of whether the formalism is valid may not be determined from the data. A graphic example of this is discussed by Raymond (1990).

4. Instructional Examples

We present three examples. The first serves to emphasize the non-uniqueness of the interpretation of emission lines in terms of plasma density and temperature structure, for a case in which the atmospheric structure is given, and traditionally made assumptions are met, based upon work of Judge, Hubený and Brown (1997). We simply ask the question, ‘can we tell if the solar TR and corona is in hydrostatic equilibrium’ (approximately constant pressure)? The second is a simple extension to the concept of a ‘filling factor’. The third serves to illustrate problems that can arise when trying to diagnose plasma conditions when the plasma evolves in response to dynamic heating, and is based on the work of Wikstøl, Judge, and Hansteen (1998).

The idea behind each example is simple: First, produce synthetic solar data; second try to diagnose physical conditions based only on the synthetic data using the standard and commonly used ‘inverse’ and ‘empirical’ techniques; and third consider how well or poorly these methods can reproduce the actual physical conditions underlying the simulations.

4.1. DIAGNOSIS OF PLASMA TEMPERATURE AND DENSITY STRUCTURE

Consider the formation of emission lines under standard ‘coronal’ conditions (these are reviewed by Judge, and Hubený, and Brown, 1997). For our purpose we can adopt the assumption that elemental abundances are constant. The line emission coefficients are then functions of electron density n and temperature T (e.g., Mason and Monsignori-Fossi, 1994; Judge, Hubený, and Brown, 1997). We can write the intensity of line i as g_i where

$$g_i = \int \int K_i(T, n) \mu(T, n) dT dn + \delta g_i, \quad (1)$$

where $\mu(T, n)$ is the source term, the emission measure differential in temperature and density, and δg_i are uncertainties.

The definition and physical meaning of $\mu(T, n)$ may seem obscure (Brown *et al.*, 1991), and we will see that it may never be determined by observations. However, we can suspend disbelief for a moment, and consider what determinations of $\mu(T, n)$ might tell us. A plasma consisting of a constant density and temperature would yield a single point, T_0, n_0 in the (T, n) plane, with $\mu(T, n) = \delta(T - T_0)\delta(n - n_0)\mu(T_0, n_0)$. A constant temperature (density) atmosphere would yield a straight vertical (horizontal) line in this plane, and a constant pressure atmosphere would correspond to a locus where $nT = \text{constant}$ in the same plane. One can imagine more complex forms of $\mu(T, n)$ as the observed volume of emitting plasma consists of separate structures, each with their own distributions of densities, temperatures in this plane.

4.1.1. Inverse approach

The ‘inverse’ approach amounts to solving for $\mu(T, n)$ given a set of measurements $\{g_i\}$. The ill-posedness of this problem is particularly bad (Brown *et al.*, 1991; Judge, Hubeny, and Brown, 1997): many solutions are compatible with the observed data. Drastic ‘additional information’ must be added *a priori* just in order to obtain a solution (i.e., regularization). Furthermore, errors in the ‘kernels’ $K_i(T, n)$, i.e., in the atomic excitation calculations, even in the absence of data noise ($\delta g_i = 0$) almost certainly preclude the possibility of solving the inverse problem (Judge, Hubeny, and Brown, 1997)³. The severity of the problems is illustrated by Judge, Hubeny, and Brown (1997) who studied the case where $\mu(T, n)$ corresponds to a solar emission measure differential in T , at a constant pressure appropriate to the quiet Sun.

Under these conditions *we cannot expect to determine $\mu(T, n)$ in cases of practical interest*, in spite of early optimism (Hubeny and Judge, 1995). We must therefore turn to see if the ‘empirical methods’ can help, recognizing that these ‘solutions’, falling short of $\mu(T, n)$, are patently not formal solutions of the inverse problem, and they must be viewed accordingly.

4.1.2. Empirical approach-line ratios

Line ratios are by far the most common way in which plasma densities are determined, with origins as far back as Menzel *et al.* (1941). The method is reviewed by Gabriel and Jordan (1971) and Mason and Monsignori-Fossi (1994). The method is explicitly ‘empirical’ in nature, aiming to determine ‘mean’ values by asking the question ‘what is the single density $\langle n \rangle_{ij}$ that is compatible with the ratio g_i/g_j ?’. In the context of Figure 1, the variable t in the figure is set to $n - \langle n \rangle_{ij}$.

To see how this method can help diagnose plasma properties under solar conditions, consider the inverse problem mentioned above (a similar example is discussed by Brown *et al.*, 1991). Proceed as follows (forward-inverse approach): (1) select a set of lines that are sensitive to n and T . (2) Choose specific $\mu(T, n)$ distributions and simulate line intensities, using standard assumptions, and derive line ratios. These are the ‘observations’. (3) Using line ratios within a given ion, set T to the ionization equilibrium temperature for the appropriate ion and solve for the single points in the (n, T) plane that are compatible with the data. (4) Make a plot of these points and compare with the ‘real’ source, $\mu(T, n)$.

Figure 2 shows results for two simple sources. The left hand panel shows that the method works well, for a simple constant pressure distribution. The points track the form of $\mu(T, n)$, and there is no doubt that the method has revealed that the adopted atmosphere is indeed a constant pressure atmosphere. The right-hand panel shows results for another source: an atmosphere with two components, each at their own pressure, contributing equally to the emission from resonance lines. The

³There is some hope that such errors might be handled through recasting the problem in terms of line ratios (McIntosh, Brown, and Judge, 1998; McIntosh, 2000). However, the ill-posedness cannot be avoided.

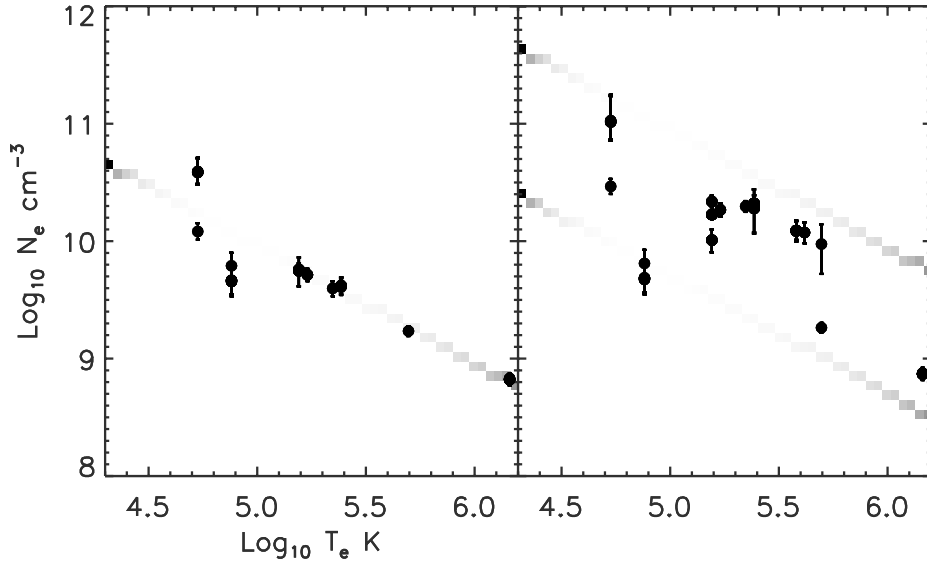


Figure 2. Density diagnostic line ratio analysis of simulated data. The points are mean densities $\langle n \rangle_{ij}$ derived using the standard techniques, using ionization equilibrium to fix the temperatures. The gray-scale images are the input emission measure function $\mu(T, n)$, corresponding to three atmospheres at different constant pressures. Although the left panel shows that the density diagnostic techniques work well for the simplest case, the mean densities found in the right panel do not exist in the source. It is also difficult to guess the form of $\mu(T, n)$ from the mean values alone.

density diagnostic approach now yields a different, more confusing picture. It is very difficult, presented with the points alone, to guess the form of the underlying source term $\mu(T, n)$.

These simple examples, which arguably represent simpler forms than are present in the actual Sun, and which *a priori* are compatible with simple assumptions (especially ionization equilibrium), illustrate clearly the non-unique interpretation of emission lines because of ill-posedness. It is important to realize that the ‘success’ of the technique that is revealed in the left hand panel of Figure 2 arises simply because $\mu(T, n)$ really is of a simple form. The failure of the method shown in the right hand panel demonstrates that this ‘success’ amounts in fact to prior knowledge of the source term: the prior information being that there is just one density at each temperature in $\mu(T, n)$. The success is judged by the fact that we knew *a posteriori* the form of $\mu(T, n)$ was simple. Ill-posedness dictates that many other forms of $\mu(T, n)$ are compatible with the data. In real astrophysical objects we do not, and may never, know the form of $\mu(T, n)$.

So much work is done using these density diagnostic line ratios that we feel obliged to re-emphasize the following points.

1. The most information one can learn in principle from a set of emission line intensities is $\mu(T, n)$.

2. The problem is so severely ill-posed that $\mu(T, n)$ may never be derived from real data.
3. The idea that ‘the simplest solution compatible with the data’, an application of Occam’s Razor, only works if you have prior knowledge of the form of $\mu(T, n)$, and if the form is particularly simple. The ‘standard line ratio technique’ is an application of this idea.

We conclude that, without other information (e.g., by applying Occam’s Razor) emission lines contain little information on density structure. In Section 5.2 we will discuss just what information might be added to mitigate these serious problems.

4.2. SUB-RESOLUTION PLASMA STRUCTURE: ‘FILLING FACTORS’

An important issue in solar (and astro-) physics concerns the nature of structure that is below spatially resolvable scales. This is not an issue that amounts to ‘details’, because much of the important physical processes (energy dissipation) must occur on unresolvable scales. One way of describing unresolved structure is to determine a filling factor. Several authors have used traditional spectroscopic techniques to determine what we will call ‘spectroscopic filling factors’ f_s based to a large degree on densities derived from the line ratio technique.

Recently, Judge (2000) investigated the meaning of f_s based upon the formalism of Almléaky, Brown, and Sweet (1989). Using ‘forward-inverse’ calculations for some *ad hoc* (but not Dirac- δ function) distributions of electron density along the line of sight, he investigated the effects of finite widths in the assumed distributions and concluded that the derived filling factors (1) systematically underestimate the true filling factor unless the plasma is truly homogeneous, (2) depend on the choice of line pairs, and (3) depart more from the true filling factor for broader distributions. Given that the form of the distribution may never be known from observations (Judge, Hubený, and Brown, 1997), this re-emphasizes the non-unique interpretation of data from unresolved plasmas. We will return to this subject in section 5.3.

4.3. DIAGNOSIS OF THE UNRESOLVED STRUCTURE OF THE TRANSITION REGION

A striking theme of the results from the TRACE and SOHO missions is that the corona and TR are dynamic in nature. Theoretical considerations also imply that heating mechanisms are expected to be dynamic (e.g., works by Nordlund, Gomez in this volume). While dynamics has been studied observationally for two or more decades (as reviewed by Mariska, 1992, ch. 6), a study that focussed on the influence of dynamics on the basic interpretation of coronal and TR data has appeared only recently (Wikstøl, Judge, and Hansteen, 1998). These authors took the following (again, forward-inverse) approach: first, build time dependent models of the corona and TR. Second, compute the emergent spectrum at each point. Third, perform suitable averages (in an attempt to mimic line of sight and instrumental

spatial and temporal integrations) of these data. Lastly, examine the synthetic data using commonly used techniques and attempt to determine the physical nature of the emitting plasma.

Wikstøl, Judge, and Hansteen (1998) chose to examine simple CTR models (ignoring cross field conduction) in which the TR is formed at the thermal interface between the corona and chromosphere, for several reasons. First, a CTR must exist on the Sun. Second, electron heat conduction has sufficient heat flux to account for all radiative losses from plasma down to $\sim 10^4$ K (e.g., Athay, 1990). Third, for a prescribed coronal temperature, the thermal structure is determined simply by a balance between heat conduction and other terms which can be accurately calculated. (This can be contrasted with the other kinds of models mentioned in Section 2.) Lastly, the CTR model has received considerable criticism, and it is important to determine the uniqueness of claims for or against such models.

Wikstøl, Judge, and Hansteen (1998) proceeded to re-examine earlier evidence cited against the dominance of CTR models. The evidence, collected in Table I, will be discussed further below. First we will show how the ‘empirical’ methods fail completely under certain conditions. Consider point 2. listed in the table, which Mariska (1992) finds the most compelling of Feldman’s (1983) arguments. In the calculations of Wikstøl, Judge, and Hansteen (1998), the TR emission at all times is formed in a thin interface between the corona and chromosphere, and yet when averaged over time the unresolved dynamics *gives the appearance* of a TR which, at the limb, is several Mm thick. This graphic demonstration of the failure of the ‘empirical’ methods, because of one or more incorrect implicit assumptions, highlights the potential danger of over-interpreting emission line data, irrespective of whether the calculations of Wikstøl, Judge, and Hansteen (1998) represent the actual Sun. A similar situation, though on firmer ground (because the forward calculations are less *ad hoc*), has arisen concerning the structure of the solar chromosphere. Carlsson and Stein (1995) have questioned the validity of stationary, semi-empirical chromospheric models (e.g., Vernazza, Avrett, and Loeser, 1981) based upon dynamic models which can reproduce important aspects of time-dependent line profiles.

On the basis of such calculations, Wikstøl, Judge, and Hansteen (1998) leveled similar criticisms at much of the other ‘empirical evidence’ for nCTR models listed in table I (in particular points 3., 4. and 9.). For example, Figure 3 illustrates that the TR and coronal lines can appear to vary very differently in time, even though the TR emission is at all times connected to the corona. In each case, the application of Occam’s Razor to search for the simplest solution that is compatible with the data yields results that are incompatible with the physical model actually used.

To conclude, Wikstøl, Judge, and Hansteen (1998) showed that the commonly used methods used to diagnose plasmas under (implicitly assumed) static conditions will fail if the TR is as dynamic as their models suggest. Incorrect conclusions can arise from applying reasonable and traditional diagnostic methods to spectral data, when unresolved dynamic evolution of the emitting plasma is important. In

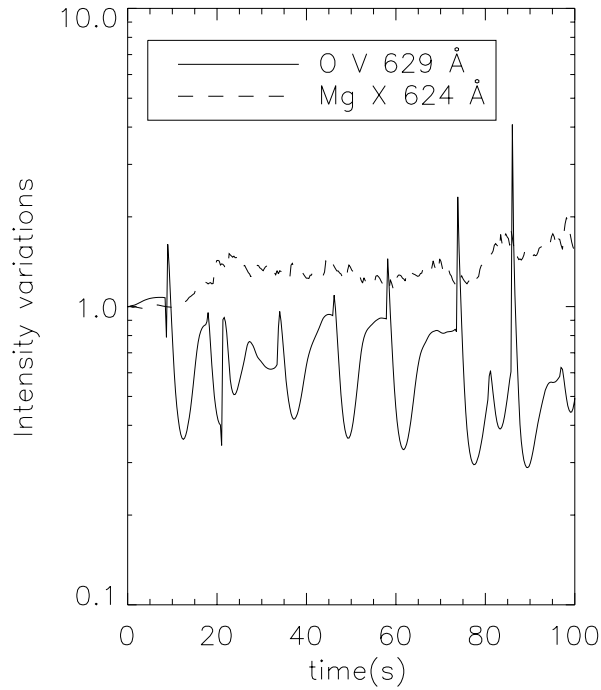


Figure 3. The figure shows total intensity, integrated over wavelength, for the O v $\lambda 629$ Å (formed near $\log T = 5.3$) and Mg x $\lambda 624$ Å ($\log T = 6$) lines as a function of time. Intensities are normalized to the initial values.

the TR, the discrepancies can be particularly dramatic, owing to the presence of steep temperature gradients. This analysis therefore casts doubt on the value of ‘empirical’ methods. More disturbingly, it reveals the possibility for a ‘collective misinterpretation’ of data in terms of a certain class of physical model. One clear example of the inability of unresolved emission line spectra to discriminate between very different classes of models was given by Raymond (1990). Without this work, one might have been led to conclude that a static picture of the TR is a good approximation, a conclusion that is not at all warranted by the data analyzed there.

5. Discussion

5.1. EMPIRICAL VS. INVERSE VS. FORWARD METHODS

The ‘empirical’, ‘inverse’ and ‘forward’ approaches share several characteristics: all require some kind of model, all are subject to non-uniqueness, and all require additional information before sensible conclusions can be drawn from the data.

TABLE I
Evidence against classical transition region models

Point	Observation	Interpretation	
		Static picture	Dynamic picture
1.	Absence of complete structures in spectroheliograms ^a	nCTR	CTR?
2.	$(d \ln I/ds)^{-1}$ (Limb intensity scale height) \gg classical TR thickness ^b	nCTR	CTR
3.	non-thermal line widths differ in TR and corona ^b	nCTR	CTR
4.	$P_e(TR)$ varies dramatically, unlike chromosphere/corona ^b	nCTR	CTR
5.	$\xi(T)$ below $T = 2 \times 10^5$ K requires heating beyond conduction ^b	nCTR	CTR
6.	Abundances differ in corona and TR ^c	nCTR	CTR?
7.	‘Threads’ in TR lines seen at the limb ^d	nCTR	CTR?
8.	Tiny spectroscopic filling factors ^e	nCTR	CTR
9.	$\partial I/\partial t$ differ in corona and TR ^f	nCTR	CTR

^aFeldman (1983, 1987), Feldman and Laming (1994), also implied by some (not all) active region C IV/171/195 data from TRACE. ^bFeldman (1983). ^cFeldman (1998). ^dSee any limb image from TRACE in H Ly α or C IV. ^eE.g., Feldman, Doschek, and Mariska (1979), Dere *et al.* (1987), Judge and Brekke (1994). ^fSee, for example, Figure 4 of Hansteen (1997).

The differences between the methods become clearer when examining specific examples, applicable to conditions believed to be present in the solar TR. First, it is clear that the ‘empirical’ techniques are very restricted sub-class of the inverse methods, for which ‘solutions’ are determined essentially by assumption. These techniques receive support primarily from application of Occam’s Razor alone, and our first example shows that they are in fact *determined* by the razor itself! We will see below that some support for these methods *might* be found in the physics of the Sun’s corona, but that this is as yet unclear. These methods therefore should not be trusted. Furthermore, if the thermal structure of the corona is more complex than assumed in the empirical approach, then demonstrable systematic errors will arise (e.g., Judge, 2000).

Second, the formal inverse method is less subjective, but suffers from the problem that the ill-posedness is so severe, with errors (both observational and theoretical) sufficiently large, that a meaningful inverse solution may not be found. Regularization of the solutions is thus required to remove large regions of solution space, and this may (or may not) be physically justified. There are cases where unphysical constraints are added just to regularize the solutions. A serious problem might be that the plasma conditions may not conform at all to posing the problem in inverse form (our third example is that of time dependent dynamic picture). Then even if the inversion can be done, it is akin to a small child ‘successfully forcing a square peg into a round hole’, and the interpretation is unclear (or incorrect). Worse yet is the possibility that the peg can be made to fit and the child continues to build, in spite of the poor foundation.

On this basis we are led, given what we know about physical conditions in the corona and TR, to consider forward methods as the safest approach. In this way, the needed ‘additional information’ is added through a set of equations, boundary conditions and assumptions, and not in the arguably more arbitrary and subjective fashions required to tackle data through the other approaches.

5.2. SOURCES OF ADDITIONAL INFORMATION

The empirical methods add information to the observational data essentially by applying Occam’s Razor to an inverse problem. The inverse methods add information through ‘regularization’ to make the problem less ill-posed. Both of these amount to subjective choice. Information should be added through studying the physics of the corona/TR.

Consider an active region ‘hot plasma loop’. What can we say about the physical conditions within the loop? From a discussion of the one-dimensional (steady-state) energy balance in coronal loops at constant pressure, Rosner, Tucker, and Vaiana (1978) derived some well-known loop scaling laws. In their picture the corona is made up of isolated mini-atmospheres consisting of loop-like flux bundles containing plasma at similar temperatures and densities. The laws relate the loop length L , pressure p and maximum (electron) temperature T as follows:

$$T \sim (p L)^{1/3} \text{ K}, \quad (2)$$

$$\xi(T) \sim \frac{p^2 L}{T} \sim \frac{T^5}{L} \text{ cm}^{-5}, \quad (3)$$

where $\xi(T)$ is the emission measure differential in $\ln(T)$ close to the maximum temperature in the loop. Thus, *given both T and L* , p and $\xi(T)$ are fixed. Consideration of the total energy lost by the corona per unit area yields

$$F_M \sim T^{7/2} L^{-1} \text{ erg cm}^{-2} \text{ s}^{-1}. \quad (4)$$

F_M is the mechanical energy flux needed to balance the total radiative losses. This relation shows that the corona acts like a thermostat: large changes in F_M at the

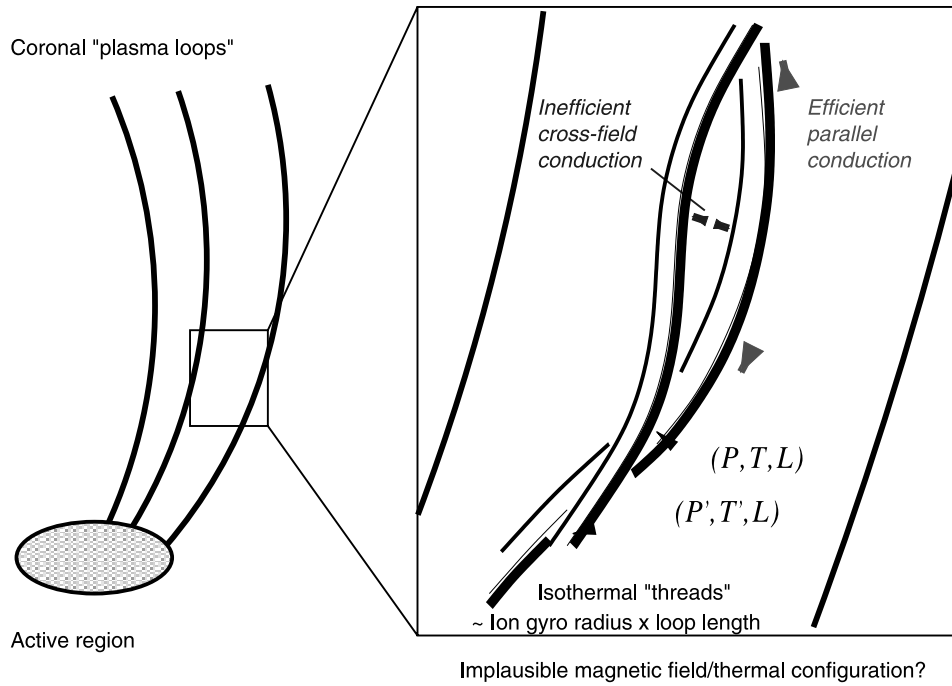


Figure 4. A cartoon illustrating the problem of the thermal structure within a given plasma loop. Is the apparently 'resolved' loop, dubbed a 'discrete coronal structure' by Rosner, Tucker, and Vaiana (1978) seen at 1 Mm spatial resolution, actually a bundle of 'micro-' or 'nano-' atmospheres, each with their own thermal properties?

base of the loop yield only small changes in T . In turn, for a given L , both p and $\xi(T)$ are fixed by a given value of F_M . Thus, simple (1D) energy balance considerations suggest that, unless F_M is very different within a given loop, *plasma loops should indeed appear to be rather homogeneous*. This argument implicitly assumes a steady state heating mechanism.

We must however ask the question, what is the likely 3D distribution of T *within* the loop? Litwin and Rosner (1993) argued that cross field transport is potentially a problem because, in spite of the thermostat noted above, dissipation of magnetic energy must occur in a tiny volume of a given loop, owing to the enormous magnetic Reynolds numbers in the corona. Classical heat conduction cannot effectively transport heat across field lines. Therefore, unless there is some other mechanism transporting heat across field lines *within* a loop, or the heating mechanism itself leads to heat transport across field lines, one would expect thermal properties to differ dramatically *within* a loop, on extremely small scales. This would lead to a picture such as shown in Figure 4. In this case⁴ a given loop contains a distribution of fine scale 'loops', each with their own values of p , T for the same length L . If

⁴Note that this kind of picture was invoked by Athay (1990) to account for TR emission via cross field conduction close to the footpoints of loops at different coronal temperatures.

such structure is present on the Sun, the simple plasma diagnostic techniques will fail.

High resolution observations show that the corona is more organized than the picture in Figure 4, where plasma ‘threads’ of different T , p co-exist within one resolvable ‘loop’, suggests. Specifically, TRACE and NIXT images reveal that large volumes of plasma in active region loops are predominantly at the same temperature to within a factor of two. For example, active regions seen in the TRACE bands at 171, 195 and 284 Å typically show loop structures in each band of several Mm apparent width, but each such loop can be physically separated in space from neighboring loops. Spatially coincident loops visible in even two of these bands are rarely observed, as Figure 4 might naively suggest. This suggests that the Sun’s corona likes to organize itself into the mini-atmospheres mentioned above, at least in a qualitative sense. It would be interesting to set some hard limits on the relative amounts of material at say 1 MK and 2 MK within certain loop structures.

While the temperature structure within a resolved plasma loop is not currently well constrained, it seems clear that there is indeed unresolved structure, simply from the well-known fact that line-widths always exceed thermal values. Furthermore, recent work has suggested small filling factors (DiMatteo *et al.*, 1999) based upon a technique that avoids the problems mentioned in Section 4.2, although the interpretation is (naturally) model dependent.

We conclude that no definitive answer can be provided as to the internal structure of plasma loops seen in the corona. It is not possible to show with confidence that T and p are uniform from current observations. Furthermore, if there are other types of structure along the line of sight (for example, low lying cool loops of Dowdy, Rabin, and Moore, 1986; Antiochos and Noci, 1986) that contribute to the observed intensities, then one cannot expect uniformity of T and p . In this situation we must be even more careful not to over-interpret data in terms of traditional plasma diagnostic techniques.

5.3. YET ANOTHER LOOK AT THE OLD PROBLEM OF TRANSITION REGION STRUCTURE

In view of the non-uniqueness issues discussed above, it is important for us to revisit the question: what is the essential structure of the solar TR? The literature currently provides evidence in support of two types of models: the CTR (thermal interface) model, including the cross-field conduction models of Athay (1990) and Ji, Song, and Hu (1996), and the other models (non-CTR) which are presumably all of ‘cool loop’ form (e.g., Antiochos and Noci, 1986), or they are transient in nature (e.g., Sturrock *et al.*, 1990; Spadaro, Lanza, and Antiochos, 1996). We will deliberately sidestep the (important) issue of energy balance in this discussion, because we wish to avoid making implicit assumptions⁵.

⁵Arguments based upon energy balance are useful once the basic structure is known. There is nothing terribly wrong with avoiding the energy balance issue at this stage, since in any case even

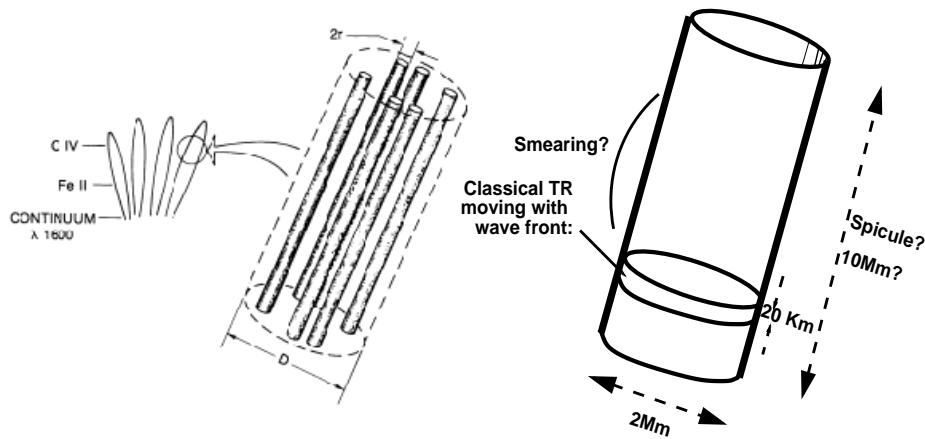


Figure 5. Two extreme views of the unresolved structure of the solar transition region compatible with the same observational data. The left picture shows the filamentary, non- classical transition region structure suggested by Dere *et al.* (1987), the right panel the dynamic CTR picture of Judge (2000).

Debate continues as to the contributions of CTR or non-CTR models to the observed solar TR emission, illustrating again the non-uniqueness issue. The CTR ‘interface’ models have found support in recent work by Gallagher *et al.* (1998), for the quiet Sun. The non-CTR class of models finds recent support from TRACE observations near active regions, which sometimes show what appears to be similar TR structures in solar images, but with quite different overlying coronal conditions (T. Berger, private communication), emphasizing some earlier work with SKYLAB data (Feldman, 1983, 1987; Feldman and Laming, 1994). Let us then try to resolve this debate and synthesize a picture of the solar TR.

Two extreme views of the unresolved geometry of the emitting plasma are presented by Dere *et al.* (1987), and by Judge (2000). These are shown (in cartoon form) in Figure 5. Dere *et al.* (1987) have argued for a highly filamentary structure based upon the (indisputable) fact that images reveal the TR to be structured into Mm length scales, and the (debatable) analysis of density sensitive emission lines. Judge (2000) has argued for a CTR model, but one in which the TR moves dynamically (and dramatically) in response to very sporadic episodes of coronal heating, based somewhat on the work of Wikstøl, Judge, and Hansteen (1998).

The argument in favor of the picture derived by Dere *et al.* (1987), is quite simple: adopting the ‘observed volume’ of the emitting plasmas (determined from high resolution images, and including the center-to-limb variations), and the spectroscopically determined density, it is clear that the volume of the plasma, if filled with material at the ‘measured density’, would emit between 100 and 10^5 times the amount of radiation compared with what is observed. The conclusion is that just a

if electron conduction redistributes energy in specific ways, this begs the question of what supplies heat to the corona.

tiny fraction of the available volume is filled with TR-emitting plasma. Given the apparent vertical lengths of the structures seen at the solar limb, the only possible resolution of these facts is to force the emission into highly filamentary strands, as shown in the left panel of Figure 5. Dere *et al.* (1987) derived path lengths Δh from the equation

$$I^{Obs} = G(T)n_e^2\Delta h, \quad (5)$$

where I^{Obs} is the observed intensity of the C IV resonance line, and n_e was derived using the density diagnostic line ratios of O IV, assuming pressure equilibrium. The derived path lengths vary from 0.1 to 10 km.

Judge (2000) presents a different interpretation of these same observational facts, arguing that dynamics moves the CTR along field lines, spatially smearing the TR radiation in that direction. He also argues that the path lengths are formally lower limits owing to the systematic effect discussed in Section 4.2. Values close to the scale height of classically heated TR models are not unreasonable. Two other facts, difficult to account for in the filamentary picture, may also fall into place in a dynamic CTR picture: Dere *et al.* noted that the highest pressure regions have the shortest path lengths, a natural consequence of CTR models. They also emphasized that it is difficult to explain why, in the filamentary picture, large areas (several Mm along the slit) share the same large-scale velocities. The CTR picture might also explain this, if the corona is (for reasons not yet understood) horizontally uniform across similar length scales.

The striking differences between these two pictures highlights the ambiguities concerned with determining the nature of spatially unresolved structures using spectroscopy, in the absence of other information. Indeed, one strong possibility is that there is some truth in both pictures! The physical conditions in the TR and corona present very challenging problems to theorists, so, unlike the photosphere (for instance), additional information cannot yet be obtained from models. Accordingly, the resolution of this problem must come from higher (sub-arcsecond) spatial resolution observations not yet possible with existing instruments.

Returning to Table I, we are left with points 1., 6., and 7. (marked CTR? in the table) as the evidence against CTR models. Point 1. was addressed by Wikstøl, Judge, and Hansteen (1998) by arguing that the physical volumes emitting coronal and TR lines differ dramatically because of the very steep dependence of temperature gradient on temperature. In essence, coronal emission lines form over very large volumes which, in quiet Sun conditions, will smear out coronal images and make them appear (even for simple magnetic geometries) very different from images of TR lines. Furthermore, recent high resolution loop images have, in active regions, shown direct evidence for the link between plasma at 10^6 and a few times 10^6 K (Berger *et al.*, 1999) suggestive of electron conduction at least to the ‘top’ of the TR, and there is some worry that the foot-points of coronal loops are masked by absorption by spatially mixed cooler material evident in $H\alpha$ images (e.g., Peres, Reale, and Golub, 1994; Berger *et al.*, 1999). Ji, Song, and Hu

(1996) also showed that, by allowing cross field conduction in different geometries near loop footpoints, radically different $\xi(T)$ functions (see their Figure 4) and hence relative coronal and TR line intensities can be produced all in the context of (cross-field) conductively heated models. We conclude that more information than just images is needed to judge connectivity. It remains to be seen if point 6. can be shown to be inconsistent with CTR models, especially since (1) the abundances are determined through questionable spectral diagnostic techniques and (2) element fractionation is expected in the presence of steep temperature gradients owing to the dominant effects of the thermal force. Point 7. shows that the Sun can produce very long filamentary structures that probably cannot be explained by a simple CTR model. This too cannot be used to argue against conductively heated models because one might naturally expect such features at the interface between adjacent coronal flux bundles, at different temperatures, as a result of cross field conduction (Athay, 1990). In any case the observations indicate that the contribution of such features, seen only at the limb, to the disk intensity, is small.

Lastly, we believe there is one ‘red herring’ that has been cleared up through studying statistically large samples of data from the SUMER instrument on SOHO. Interesting models of Antiochos (1984) and Spadaro, Lanza, and Antiochos (1996), invoking radiative transfer processes, adopted the ‘observational result’ that TR lines are red-shifted everywhere on the disk, with *no* $\cos \vartheta$ dependence, based primarily upon the active region observations of Feldman, Cohen, and Doschek (1982), and arguments presented by Feldman (1983). Peter and Judge (1999) have demonstrated that representative UV lines spanning temperatures from 10^4 to 10^6 K in the quiet Sun obey statistically the $\cos \vartheta$ dependence expected for optically thin emission. Thus, there is no need to go to unusual lengths to account for such behavior, at least for the quiet Sun.

6. Conclusions

Information is not in the data alone- we require additional information to determine reliably properties of the emitting plasmas. There is no deep philosophical difference between ‘inverse’ and ‘forward’ approaches. In fact, models/assumptions are needed for both, but inverse methods require more restrictive assumptions (the spectrum formation must fit within the ‘inverse problem’ formalism). Popular ‘spectral diagnostic techniques’, classified here under ‘empirical approaches’, are shown to be (drastically) simplified applications of the ‘inverse’ approach, with associated problems. Occam’s Razor applied to spectral diagnostic methods is appealing but worrisome, given that the entire ‘solution’ is determined almost completely by the ‘razor’ alone. Forward models are tractable only in limited cases. Simple examples (forward-inverse calculations) reveal potential pitfalls with the inverse approaches, and so the additional data must instead be added through physical constraints using forward models. However, our understanding of the

physics of the corona/TR is very incomplete. Although the 'loop scaling laws' of Rosner, Tucker, and Vaiana (1978) (and others) would suggest relatively uniform loop properties (given the weak dependence of loop temperature on mechanical energy flux), and current observations also suggest that the corona is quite well organized on resolvable (Mm) scales, significant unresolved structure is expected on the basis of the physics of energy dissipation and transport. Thus one cannot expect simple plasma diagnostic techniques to be compatible with the real corona, where finite width distributions of plasma temperature, density and velocity are expected (and observed) in observable volumes.

Under these conditions, the only safe approach appears to be to use forward models, with all their restrictions and non-uniqueness, as a guide to the interpretation of data from the corona and TR. In this way we hope to avoid the issue of a 'collective mis-interpretation' of the data acquired at such great expense. In this sense we see the corona and TR in a similar light as Carlsson and Stein (1995) view the chromosphere.

Acknowledgements

We are grateful to Leon Golub and Tom Berger for useful discussions.

References

- Almleaky, Y. M., Brown, J. C., and Sweet, P. A.: 1989, *Astron. Astrophys.* **224**, 328.
 Anderson-Huang, L. S.: 1998, *Space Sci. Rev.* **85**, 203.
 Antiochos, S. K.: 1984, *Astrophys. J.* **280**, 416.
 Antiochos, S. K. and Noci, G.: 1986, *Astrophys. J.* **301**, 440.
 Athay, R.: 1990, *Astrophys. J.* **362**, 364.
 Berger, T. E., De Pontieu, B., Schrijver, C. J., and Title, A. M.: 1999, *Astrophys. J.* **519**, L97.
 Brage, T., Judge, P. G., and Brekke, P.: 1996, *Astrophys. J.* **464**, 1030.
 Brown, J. C., Dwivedi, B. N., Almleaky, Y. M., and Sweet, P. A.: 1991, *Astron. Astrophys.* **249**, 277.
 Cally, P. S.: 1990, *Astrophys. J.* **355**, 693.
 Cally, P. S. and Robb, T. D.: 1991, *Astrophys. J.* **372**, 329.
 Carlsson, M. and Stein, R. F.: 1995, *Astrophys. J.* **440**, L29.
 Craig, I. J. D. and Brown, J. C.: 1986, *Inverse Problems in Astronomy*, Hilger, Bristol.
 Dere, K. P., Bartoe, J.-D. F., Brueckner, G. E., Cook, J. W., and Socker, D. G.: 1987, *Solar Phys.* **114**, 223.
 Di Matteo, V., Reale, F., Peres, G., and Golub, L.: 1999, *Astron. Astrophys.* **342**, 563.
 Dowdy, J. F. J., Rabin, D., and Moore, R. L.: 1986, *Solar Phys.* **105**, 35.
 Feldman, U.: 1983, *Astrophys. J.* **275**, 367.
 Feldman, U.: 1987, *Astrophys. J.* **320**, 426.
 Feldman, U.: 1998, *Astrophys. J.* **507**, 974.
 Feldman, U. and Laming, J. M.: 1994, *Astrophys. J.* **434**, 370.
 Feldman, U., Cohen, L., and Doschek, G.: 1982, *Astrophys. J.* **255**, 325.
 Feldman, U., Doschek, G. A., and Mariska, J. T.: 1979, *Astrophys. J.* **229**, 369.
 Gabriel, A.: 1976, *Phil Trans. Royal Soc. London* **281**, 339.

- Gabriel, A. H. and Jordan, C.: 1971, *Case Studies in Atomic Collision Physics*, Ch. 4, North-Holland, pp. 210–291.
- Gallagher, P. T., Phillips, K. J. H., Harra-Murnion, L. K., and Keenan, F. P.: 1998, *Astron. Astrophys.* **335**, 733.
- Hansteen, V.: 1997, in A. Wilson (ed.), *The Fifth SOHO Workshop, The Corona and Solar Wind Near Minimum Activity*, ESA SP-404, ESTEC, Noordwijk, the Netherlands, p. 45.
- Harrison, R. A. and Thompson, A. M.: 1991, Intensity Integral Inversion Techniques: a Study in Preparation for the SOHO Mission, Technical Report RAL-91-092, Rutherford Appleton Laboratory.
- Hubený, V. and Judge, P. G.: 1995, *Astron. J.* **448**, L61.
- Ji, H. S., Song, M. T., and Hu, F. M.: 1996, *Astrophys. J.* **464**, 1012.
- Judge, P. G.: 2000, *Astrophys. J.*, in press.
- Judge, P. G. and Brekke, P.: 1994, in K. S. Balasubramaniam and G. Simon (eds.), *The 14th International Summer Workshop: Solar Active Region Evolution – Comparing Models with Observations*, Astronomical Society of the Pacific, San Francisco CA, p. 321.
- Judge, P. G., Hubený, V., and Brown, J. C.: 1997, *Astrophys. J.* **475**, 275.
- Judge, P. G., Woods, T. N., Brekke, P., and Rottman, G. J.: 1995, *Astrophys. J.* **455**, L85.
- Litwin, C. and Rosner, R.: 1993, *Astrophys. J.* **412**, 375.
- Mariska, J. T.: 1992, *The Solar Transition Region*, Cambridge University Press, Cambridge.
- Mason, H. E. and Monsignori-Fossi, B. C.: 1994, *Astron. Astrophys. Rev.* **6**, 123.
- McIntosh, S., Brown, J. C., and Judge, P. G.: 1998, *Astron. Astrophys.* **333**, 333.
- McIntosh, S. W.: 2000, *Astrophys. J.*, in press.
- Menzel, D. H., Aller, L. H., and Hebb, M. H.: 1941, *Astrophys. J.* **93**, 230.
- Parker, E. N.: 1988, *Astrophys. J.* **330**, 474.
- Peres, G., Reale, F., and Golub, L.: 1994, *Astrophys. J.* **422**, 412.
- Peter, H. and Judge, P. G.: 1999, *Astrophys. J.* **522**, 1148.
- Rabin, D. and Moore, R.: 1984, *Astrophys. J.* **285**, 359.
- Raymond, J. C.: 1990, *Astrophys. J.* **365**, 387.
- Rosner, R., Tucker, W. H., and Vaiana, G. S.: 1978, *Astrophys. J.* **220**, 643.
- Roumeliotis, G.: 1991, *Astrophys. J.* **379**, 392.
- Spadaro, D., Lanza, A. F., and Antiochos, S. K.: 1996, *Astrophys. J.* **462**, 1011.
- Steiner, O., Grossmann-Doerth, U., Knoelker, M., and Schuessler, M.: 1998, *Astrophys. J.* **495**, 468.
- Sturrock, P. A., Dixon, W. W., Klimchuk, J. A., and Antiochos, S. K.: 1990, *Astrophys. J.* **356**, L31.
- Vernazza, J. E., Avrett, E. H., and Loeser, R.: 1981, *Astrophys. J. Suppl.* **45**, 635.
- Wikstøl, Ø., Judge, P. G., and Hansteen, V.: 1998, *Astrophys. J.* **501**, 895.